

# Using Application-Specific Performance Models to Inform Dynamic Scheduling

Jeffrey S. Vetter

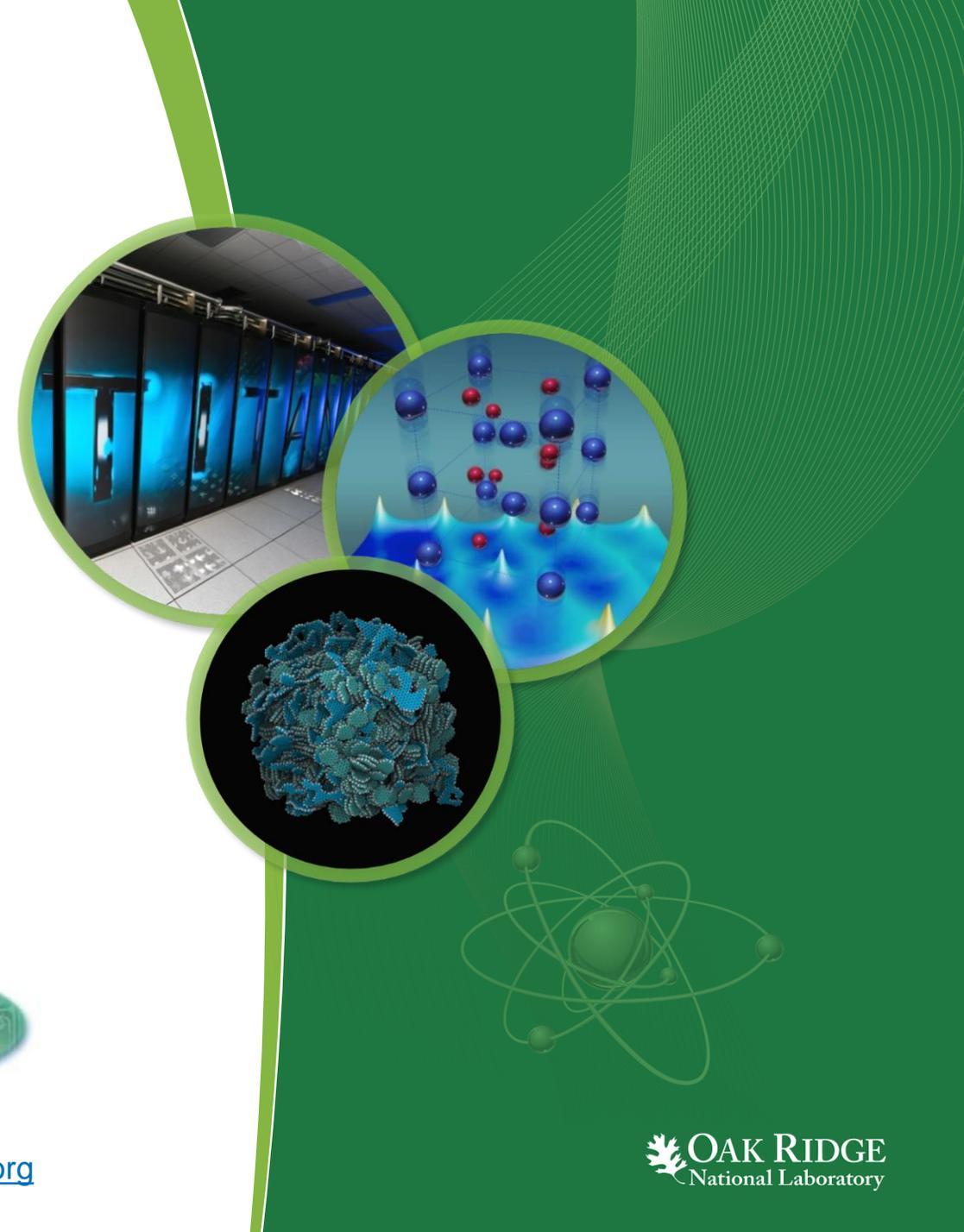
Seyong Lee, Jeremy Meredith,  
and many collaborators

*Presented to*  
**International Workshop on  
Runtime and Operating Systems for  
Supercomputers ROSS 2016**

1 June 2016  
Kyoto



<http://ft.ornl.gov> [vetter@computer.org](mailto:vetter@computer.org)



# Executive Summary

- Both architectures and applications are growing more complex
  - Trends dictate that this will get worse, not better
  - This complexity creates irregularity in computation, communication, and data movement
- Dynamic resource management is one way to help manage this irregularity
  - Need accurate policies to guide resource decisions
  - Examples: greedy work stealing, algorithmic, historical, cost models, application specific, etc
- Posit that we can use application-specific performance models to inform scheduling decisions
  - Aspen performance modeling language helps create models
  - Two recent experiments
    - GPU offload
    - *Distributed scientific workflows*

# Trends toward Exascale

# Exascale architecture targets circa 2009

## 2009 Exascale Challenges Workshop in San Diego

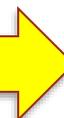
**Attendees envisioned two possible architectural swim lanes:**

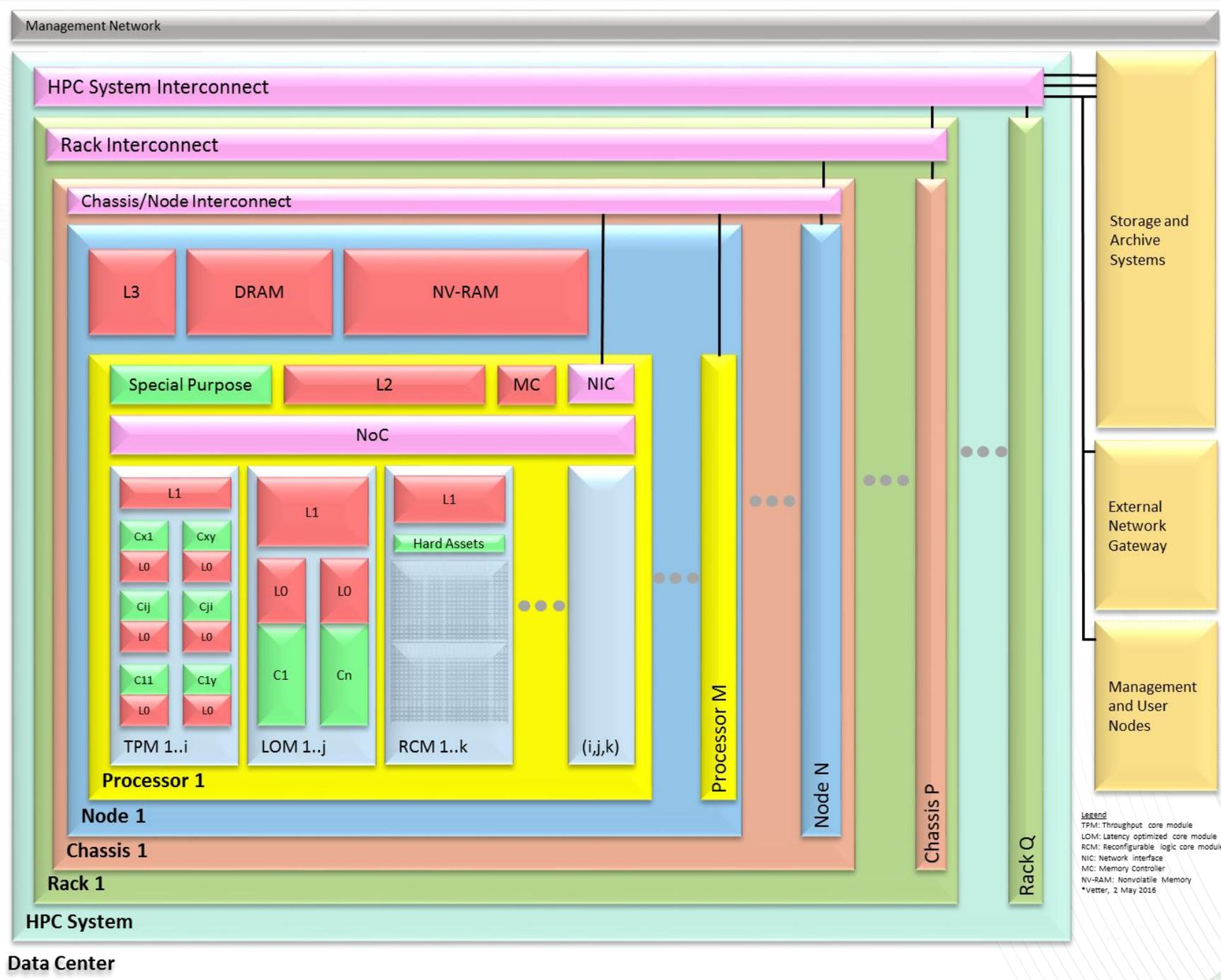
1. Homogeneous many-core thin-node system
2. Heterogeneous (accelerator + CPU) fat-node system

System attributes	2009	"Pre-Exascale"		"Exascale"	
System peak	2 PF	100-200 PF/s		1 Exaflop/s	
Power	6 MW	15 MW		20 MW	
System memory	0.3 PB	5 PB		32-64 PB	
Storage	15 PB	150 PB		500 PB	
Node performance	125 GF	0.5 TF	7 TF	1 TF	10 TF
Node memory BW	25 GB/s	0.1 TB/s	1 TB/s	0.4 TB/s	4 TB/s
Node concurrency	12	O(100)	O(1,000)	O(1,000)	O(10,000)
System size (nodes)	18,700	500,000	50,000	1,000,000	100,000
Node interconnect BW	1.5 GB/s	150 GB/s	1 TB/s	250 GB/s	2 TB/s
IO Bandwidth	0.2 TB/s	10 TB/s		30-60 TB/s	
MTTI	day	O(1 day)		O(0.1 day)	

# Contemporary ASCR Computing At a Glance

System attributes	NERSC Now	OLCF Now	ALCF Now	NERSC Upgrade	OLCF Upgrade	ALCF Upgrades	
Planned Installation	<b>Edison</b>	<b>TITAN</b>	<b>MIRA</b>	<b>Cori 2016</b>	<b>Summit 2017-2018</b>	<b>Theta 2016</b>	<b>Aurora 2018-2019</b>
System peak (PF)	2.6	27	10	> 30	150	>8.5	180
Peak Power (MW)	2	9	4.8	< 3.7	10	1.7	13
Total system memory	357 TB	710TB	768TB	~1 PB DDR4 + High Bandwidth Memory (HBM)+1.5PB persistent memory	> 1.74 PB DDR4 + HBM + 2.8 PB persistent memory	>480 TB DDR4 + High Bandwidth Memory (HBM)	> 7 PB High Bandwidth On-Package Memory Local Memory and Persistent Memory
Node performance (TF)	0.460	1.452	0.204	> 3	> 40	> 3	> 17 times Mira
Node processors	Intel Ivy Bridge	AMD Opteron Nvidia Kepler	64-bit PowerPC A2	Intel Knights Landing many core CPUs Intel Haswell CPU in data partition	Multiple IBM Power9 CPUs & multiple Nvidia Voltas GPUS	Intel Knights Landing Xeon Phi many core CPUs	Knights Hill Xeon Phi many core CPUs
System size (nodes)	5,600 nodes	18,688 nodes	49,152	9,300 nodes 1,900 nodes in data partition	~3,500 nodes	>2,500 nodes	>50,000 nodes
System Interconnect	Aries	Gemini	5D Torus	Aries	Dual Rail EDR-IB	Aries	2 <sup>nd</sup> Generation Intel Omni-Path Architecture
File System	7.6 PB 168 GB/s, Lustre®	32 PB 1 TB/s, Lustre®	26 PB 300 GB/s GPFS™	28 PB 744 GB/s Lustre®	120 PB 1 TB/s GPFS™	10PB, 210 GB/s Lustre initial	150 PB 1 TB/s Lustre®

Complexity  $\propto T$  



# Complexity is the next major challenge!

- “Exciting” times in computer architecture
  - Heterogeneous cores
  - Multimode memory systems
  - Fused memory systems
  - I/O architectures
  - Error correction
  - Changing system balance
- Uncertainty, Ambiguity
  - How do we design future systems so that they are faster than current systems on mission applications?
    - Entirely possible that the new system will be slower than the old system!
  - How do we provide some level of performance portability for applications teams?
  - How do we understand reliability and performance problems?
- Managing complexity is our main challenge!

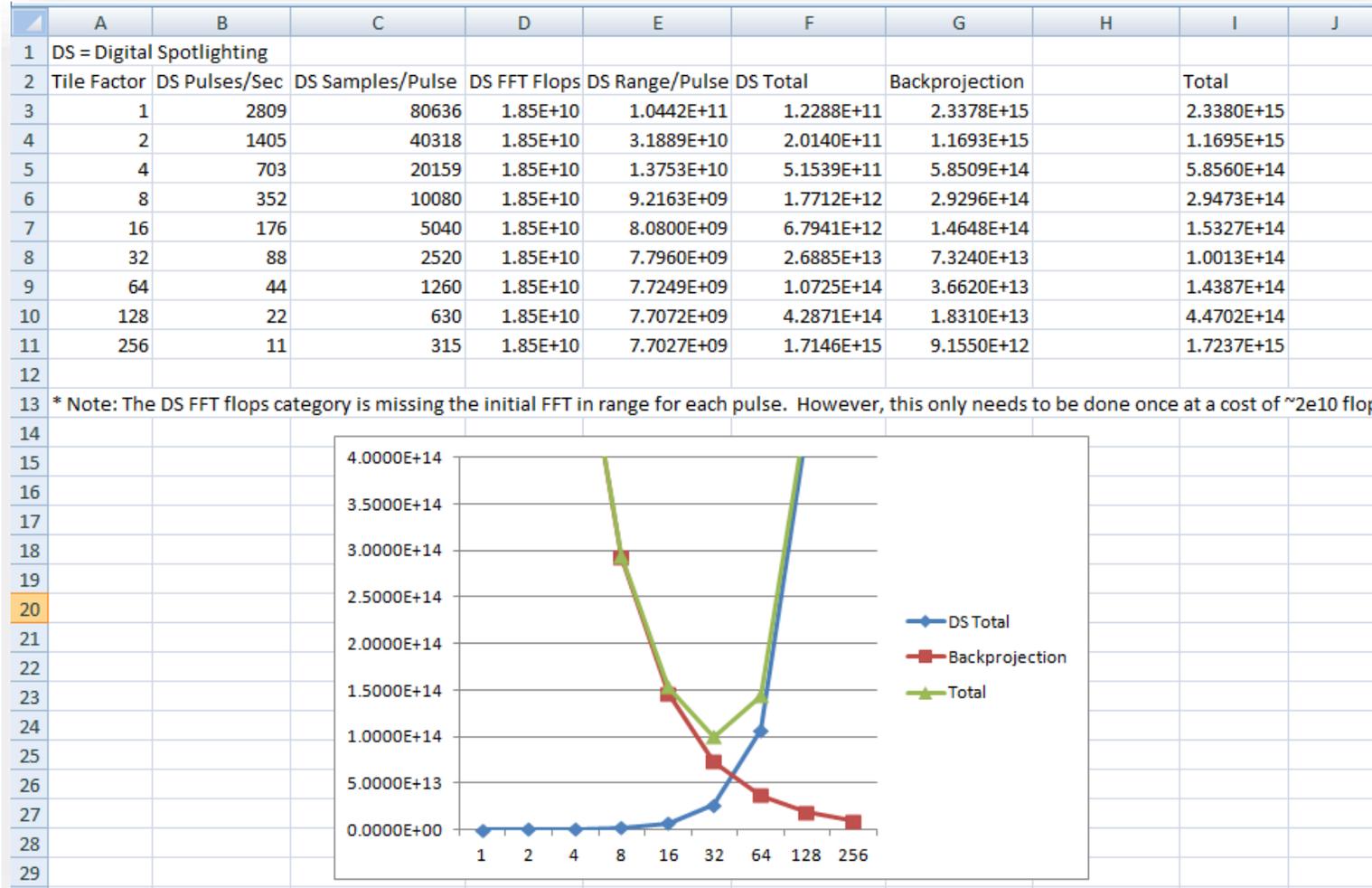
# Performance Prediction with Aspen

## Example Ad Hoc Model: Latex Equations

the communication and computation of two sheets. The expression we get for the runtime is

$$T = 2 \left[ t_c \frac{n}{p} n \log_2 n + (p - 1)o + (p - 2)g + \frac{n}{p}nG \right. \\ \left. + L + \left( \frac{n}{p} - 1 \right) \max \left\{ (p - 1)o + t_c \frac{n}{p} n \log_2 n, \right. \right. \\ \left. \left. (p - 1)g + \frac{n}{p}nG + L \right\} \right] + t_c \frac{n^2}{p^2} n \log_2 n$$

# Example: Ad-Hoc Excel Files



# Prediction Techniques Ranked

	Speed	Ease	Flexibility	Accuracy	Scalability
Ad-hoc Analytical Models	1	3	2	4	1
Structured Analytical Models	1	2	1	4	1
<i>Aspen</i>	1	1	1	4	1
Simulation – Functional	3	2	2	3	3
Simulation – Cycle Accurate	4	2	2	2	4
Hardware Emulation (FPGA)	3	3	3	2	3
Similar hardware measurement	2	1	4	2	2
Node Prototype	2	1	4	1	4
Prototype at Scale	2	1	4	1	2
Final System	-	-	-	-	-

# Aspen: Abstract Scalable Performance Engineering Notation

## Model Creation

- Static analysis via compiler, tools
- Empirical, Historical
- Manual (for future applications)



## Representation in Aspen

- Modular
- Sharable
- Composable
- Reflects graph structure

## Model Uses

- Interactive tools for graphs, queries
- Design space exploration
- Workload Generation
- Feedback to Runtime Systems

*E.g., MD, UHPC CP 1, Lulesh, 3D FFT, CoMD, VPFFT, ...*

## Source code

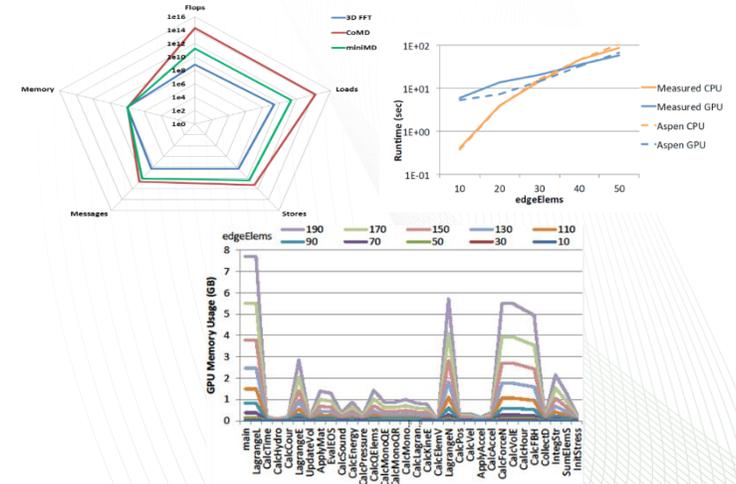
```

2324 static inline
2325 void CalcMonotonicQGradientsForElems(Index_t p_nodelist[T_NUMELEM8],
2326 Real_t p_x[T_NUMNODE], Real_t p_y[T_NUMNODE], Real_t p_z[T_NUMNODE],
2327 Real_t p_xd[T_NUMNODE], Real_t p_yd[T_NUMNODE], Real_t p_zd[T_NUMNODE],
2328 Real_t p_volo[T_NUMELEM], Real_t p_vnew[T_NUMELEM],
2329 Real_t p_delx_zeta[T_NUMELEM], Real_t p_delv_zeta[T_NUMELEM],
2330 Real_t p_delx_xi[T_NUMELEM], Real_t p_delv_xi[T_NUMELEM],
2331 Real_t p_delx_eta[T_NUMELEM], Real_t p_delv_eta[T_NUMELEM])
2332 {
2333     Index_t i;
2334     Index_t numElem = m_numElem;
2335     #pragma acc parallel loop independent present(p_vnew, p_nodelist, p_x, p_y, p_z, p_xd, \
2336 p_yd, p_zd, p_volo, p_delx_xi, p_delx_eta, p_delx_zeta, p_delv_xi, p_delv_eta, \
2337 p_delv_zeta)
2338     for (i = 0 ; i < numElem ; ++i ) {
2339         const Real_t ptiny = 1.e-36 ;
2340         Real_t ax, ay, az ;
2341         Real_t dxv, dyv, dzv ;
2342
2343         const Index_t *elemToNode = &p_nodelist[8*i];
2344         Index_t n0 = elemToNode[0] ;
2345         Index_t n1 = elemToNode[1] ;
2346         Index_t n2 = elemToNode[2] ;
2347         Index_t n3 = elemToNode[3] ;
2348         Index_t n4 = elemToNode[4] ;
2349         Index_t n5 = elemToNode[5] ;
2350         Index_t n6 = elemToNode[6] ;
2351         Index_t n7 = elemToNode[7] ;
2352
2353         Real_t x0 = p_x[n0] ;
    
```

## Aspen code

```

147 kernel CalcMonotonicQGradients {
148     execute [numElems]
149     {
150         loads [8 * indexWordSize] from nodelist
151         // Load and cache position and velocity.
152         loads/caching [8 * wordSize] from x
153         loads/caching [8 * wordSize] from y
154         loads/caching [8 * wordSize] from z
155
156         loads/caching [8 * wordSize] from xvel
157         loads/caching [8 * wordSize] from yvel
158         loads/caching [8 * wordSize] from zvel
159
160         loads [wordSize] from volo
161         loads [wordSize] from vnew
162         // dx, dy, etc.
163         flops [90] as dp, simd
164         // delvx delvx
165         flops [9 + 8 + 3 + 30 + 5] as dp, simd
166         stores [wordSize] to delv_xeta
167         // delxi delxi
168         flops [9 + 8 + 3 + 30 + 5] as dp, simd
169         stores [wordSize] to delx_xi
170         // delvj and delvj
171         flops [9 + 8 + 3 + 30 + 5] as dp, simd
172         stores [wordSize] to delv_eta
173     }
174 }
    
```



# Creating an Aspen Model

# Manual Example of LULESH

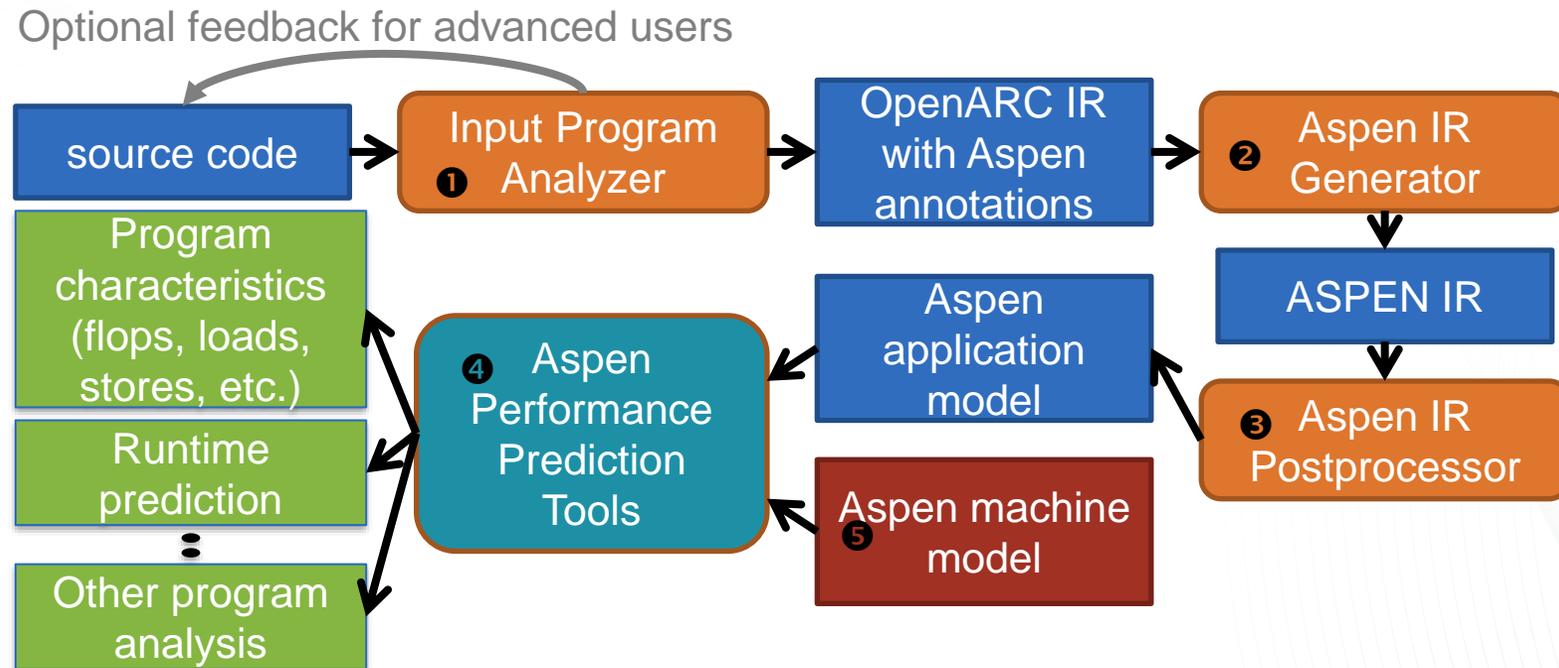
```
branch: master aspen / models / lulesh / lulesh.aspen
jsmeredith on Sep 20, 2013 adding models
1 contributor

336 lines (288 sloc) | 9.213 kb
Raw Blame History

1 //
2 // lulesh.aspen
3 //
4 // An ASPEN application model for the LULESH 1.01 challenge problem. Based
5 // on the CUDA version of the source code found at:
6 // https://computation.llnl.gov/casc/ShockHydro/
7 //
8 param nTimeSteps = 1495
9
10 // Information about domain
11 param edgeElems = 45
12 param edgeNodes = edgeElems + 1
13
14 param numElems = edgeElems^3
15 param numNodes = edgeNodes^3
16
17 // Double precision
18 param wordSize = 8
19
20 // Element data
21 data mNodeList as Array(numElems, wordSize)
22 data mMatElemList as Array(numElems, wordSize)
23 data mNodeList as Array(8 * numElems, wordSize) // 8 nodes per element
24 data mXim as Array(numElems, wordSize)
25 data mXip as Array(numElems, wordSize)
26 data mletam as Array(numElems, wordSize)
27 data mletap as Array(numElems, wordSize)
28 data mzetam as Array(numElems, wordSize)
29 data mzetap as Array(numElems, wordSize)
30 data melemBC as Array(numElems, wordSize)
31 data mE as Array(numElems, wordSize)
32 data mP as Array(numElems, wordSize)
```

```
147 kernel CalcMonotonicQGradients {
148   execute [numElems]
149   {
150     loads [8 * indexWordSize] from nodelist
151     // Load and cache position and velocity.
152     loads/caching [8 * wordSize] from x
153     loads/caching [8 * wordSize] from y
154     loads/caching [8 * wordSize] from z
155
156     loads/caching [8 * wordSize] from xvel
157     loads/caching [8 * wordSize] from yvel
158     loads/caching [8 * wordSize] from zvel
159
160     loads [wordSize] from volo
161     loads [wordSize] from vnew
162     // dx, dy, etc.
163     flops [90] as dp, simd
164     // delvk delxk
165     flops [9 + 8 + 3 + 30 + 5] as dp, simd
166     stores [wordSize] to delv_xeta
167     // delxi delvi
168     flops [9 + 8 + 3 + 30 + 5] as dp, simd
169     stores [wordSize] to delx_xi
170     // delxj and delvj
171     flops [9 + 8 + 3 + 30 + 5] as dp, simd
172     stores [wordSize] to delv_eta
173   }
174 }
```

# COMPASS System Overview



# MM example generated from COMPASS

```
1 int N = 1024;
2 void matmul(float *a, float *b, float *c){ int i, j, k ;
3 #pragma acc kernels loop gang copyout(a[0:(N*N)]) \
4 copyin(b[0:(N*N)],c[0:(N*N)])
5   for (i=0; i<N; i++){
6     #pragma acc loop worker
7     for (j=0; j<N; j++) { float sum = 0.0 ;
8       for (k=0; k<N; k++) {sum+=b[i*N+k]*c[k*N+j];}
9       a[i*N+j] = sum; }
10    } //end of i loop
11  } //end of matmul()
12 int main() {
13   int i; float *A = (float*) malloc(N*N*sizeof(float));
14   float *B = (float*) malloc(N*N*sizeof(float));
15   float *C = (float*) malloc(N*N*sizeof(float));
16   for (i = 0; i < N*N; i++)
17     { A[i] = 0.0F; B[i] = (float) i; C[i] = 1.0F; }
18 #pragma aspen modelregion label(MM)
19   matmul(A,B,C);
20   free(A); free(B); free(C); return 0;
21 } //end of main()
```

```
1 model MM {
2   param floatS = 4; param N = 1024
3   data A as Array((N*N), floatS)
4   data B as Array((N*N), floatS)
5   data C as Array((N*N), floatS)
6   kernel matmul {
7     execute matmul2_intracommIN
8     { intracomm [floatS*(N*N)] to C as copyin
9       intracomm [floatS*(N*N)] to B as copyin }
10    map matmul2 [N] {
11      map matmul3 [N] {
12        iterate [N] {
13          execute matmul5
14          { loads [floatS] from B as stride(1)
15            loads [floatS] from C; flops [2] as sp, simd }
16          } //end of iterate
17          execute matmul6 { stores [floatS] to A as stride(1) }
18          } // end of map matmul3
19        } //end of map matmul2
20      execute matmul2_intracommOUT
21      { intracomm [floatS*(N*N)] to A as copyout }
22    } //end of kernel matmul
23  kernel main { matmul() }
24 } //end of model MM
```

# Example: LULESH (10% of 1 kernel)

```
kernel IntegrateStressForElems
{
  execute [numElem_CalcVolumeForceForElems]
  {
    loads [(1*aspen_param_int)*8] from elemNodes as stride(1)
    loads [(1*aspen_param_double)*8] from m_x
    loads [(1*aspen_param_double)*8] from m_y
    loads [(1*aspen_param_double)*8] from m_z
    loads [(1*aspen_param_double)] from determ as stride(1)
    flops [8] as dp, simd
    flops [3] as dp, simd
    stores [(1*aspen_param_double)] as stride(0)
    flops [2] as dp, simd
    stores [(1*aspen_param_double)] as stride(0)
    flops [2] as dp, simd
    stores [(1*aspen_param_double)] as stride(0)
    flops [2] as dp, simd
    loads [(1*aspen_param_double)] as stride(0)
    stores [(1*aspen_param_double)] as stride(0)
    loads [(1*aspen_param_double)] as stride(0)
    stores [(1*aspen_param_double)] as stride(0)
    loads [(1*aspen_param_double)] as stride(0)
    .....
  }
}
```

- Input LULESH program:  
3700 lines of C codes
- Output Aspen model:  
2300 lines of Aspen codes

# Model Validation

	FLOPS	LOADS	STORES
<b>MATMUL</b>	15%	<1%	1%
<b>LAPLACE2D</b>	7%	0%	<1%
<b>SRAD</b>	17%	0%	0%
<b>JACOBI</b>	6%	<1%	<1%
<b>KMEANS</b>	0%	0%	8%
<b>LUD</b>	5%	0%	2%
<b>BFS</b>	<1%	11%	0%
<b>HOTSPOT</b>	0%	0%	0%
<b>LULESH</b>	0%	0%	0%

0% means that prediction fell between measurements from optimized and unoptimized runs of the code.

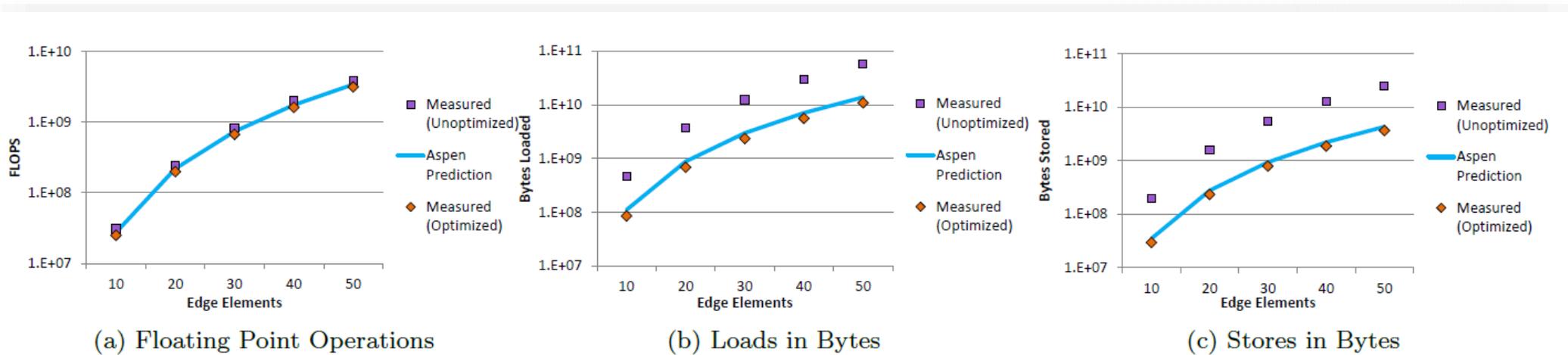
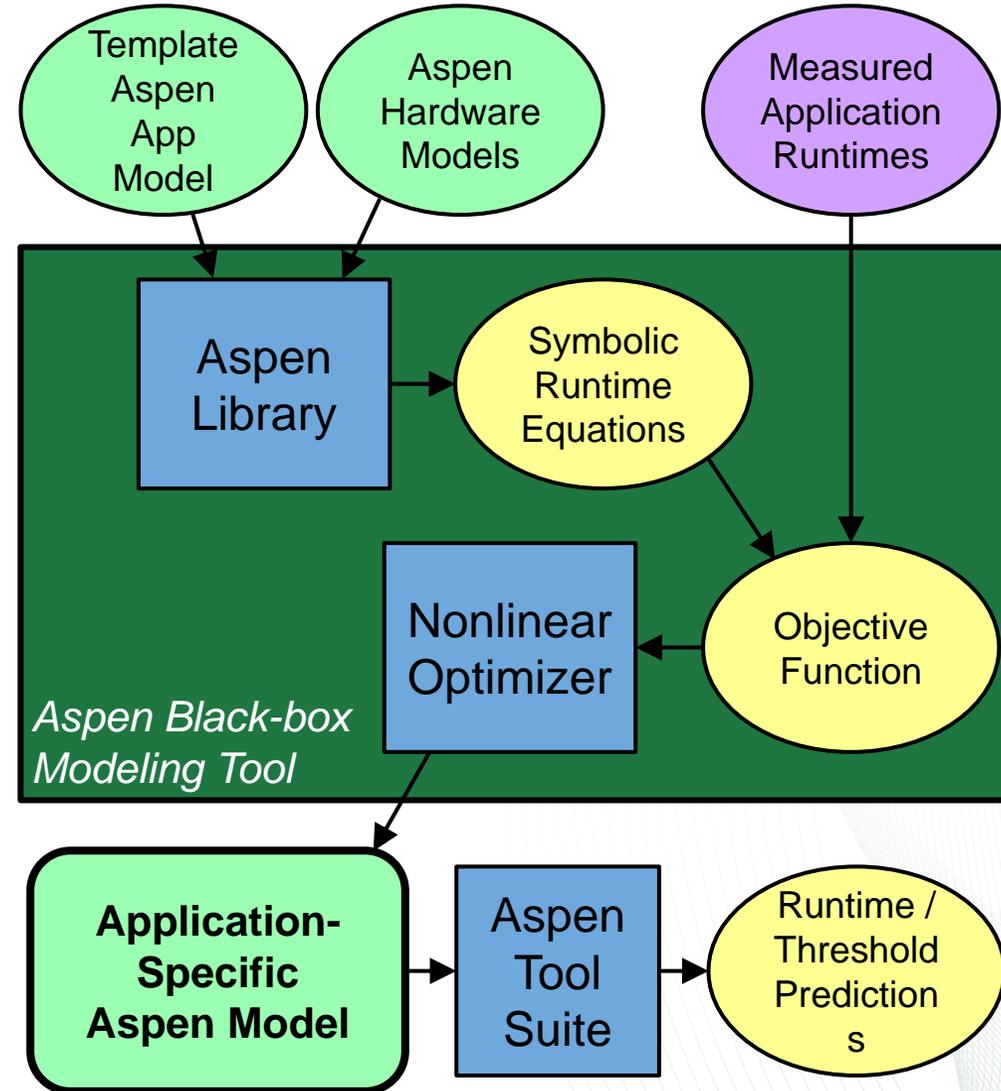


Figure 2: Predicted Resource Usage of LULESH versus Measured (with and without compiler optimization)

# Black Box Analytical Modeling

- In some cases, we do not have access to a white box Aspen performance model
- Using input vector and empirical results, we can develop Aspen Black Box model
- User provides
  - measured runtimes with app/machine parameters
    - e.g. nAtoms, nCores
  - template Aspen model with
    - application parameters
    - unknowns to solve for
  - new machine models (if necessary)
- Modeling tool
  - Generates symbolic predictions
  - Combines with measurements to generate objective function
  - Solves for unknowns in template
  - Output: completed app model usable for predictive behavior



# Black Box Modeling Example

## MD template model

```
model NAMD_Template {  
  // application parameters  
  // (defined in the input file)  
  param nAtoms = 1e6  
  param nTimeSteps = 100  
  // solve for these parameters  
  // (within the given ranges)  
  param c = 1 in 1 .. 1e18  
  param d = 1 in 1 .. 1e18  
  // application behavior:  
  // execution and control flow  
  kernel main  
  {  
    iterate [nTimeSteps] {  
      execute {  
        loads [c * nAtoms^2]  
        flops [d * nAtoms]  
      }  
    }  
  }  
}
```

## CSV data file with parameters and runtimes

nAtoms	nTimeSteps	nCores	machine	runtime
1e6	100	144	exogeni	384.2
1e6	100	144	hopper	340.1
1e6	150	144	hopper	482.9

## Concrete NAMD model

```
model NAMD_Equilibrate {  
  // NAMD input parameters  
  param nAtoms = 1e6  
  param nTimeSteps = 100  
  // calculation-specific constants  
  param c = 402.1  
  param d = 10.95  
  // NAMD application behavior  
  kernel main  
  {  
    iterate [nTimeSteps] {  
      execute {  
        loads [c * nAtoms^2]  
        flops [d * nAtoms]  
      }  
    }  
  }  
}
```

- nAtoms and nTimeSteps defined in template application model and CSV input data
- nCores defined in machine models and CSV input data
- solves for c and d, filling out a concrete application model for that problem
- new predictions can still vary nAtoms, nTimeSteps, and nCores

# Using an Aspen Model

# Aspen: Abstract Scalable Performance Engineering Notation

## Model Creation

- Static analysis via compiler, tools
- Empirical, Historical
- Manual (for future applications)



## Representation in Aspen

- Modular
- Sharable
- Composable
- Reflects graph structure

## Model Uses

- Interactive tools for graphs, queries
- Design space exploration
- Workload Generation
- Feedback to Runtime Systems

E.g., MD, UHPC CP 1, Lulesh, 3D FFT, CoMD, VPFFT, ...

## Source code

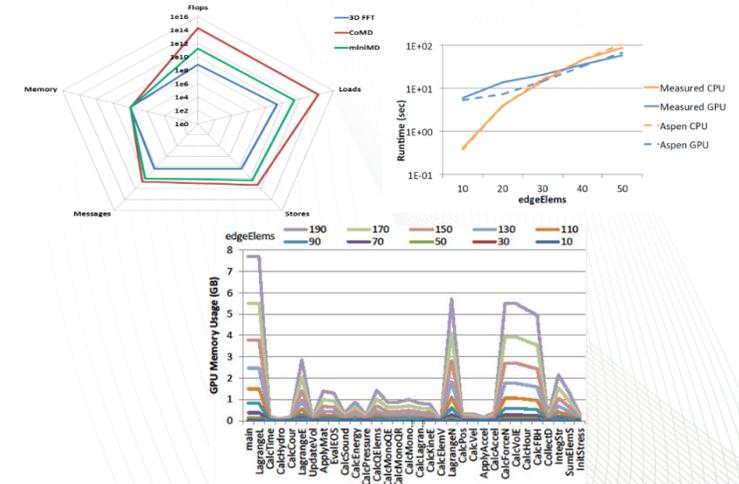
```

2324 static inline
2325 void CalcMonotonicQGradientsForElems(Index_t p_nodelist[T_NUMELEM8],
2326 Real_t p_x[T_NUMNODE], Real_t p_y[T_NUMNODE], Real_t p_z[T_NUMNODE],
2327 Real_t p_xd[T_NUMNODE], Real_t p_yd[T_NUMNODE], Real_t p_zd[T_NUMNODE],
2328 Real_t p_volo[T_NUMELEM], Real_t p_vnew[T_NUMELEM],
2329 Real_t p_delx_zeta[T_NUMELEM], Real_t p_delv_zeta[T_NUMELEM],
2330 Real_t p_delx_xi[T_NUMELEM], Real_t p_delv_xi[T_NUMELEM],
2331 Real_t p_delx_eta[T_NUMELEM], Real_t p_delv_eta[T_NUMELEM])
2332 {
2333     Index_t i;
2334     Index_t numElem = m_numElem;
2335     #pragma acc parallel loop independent present(p_vnew, p_nodelist, p_x, p_y, p_z, p_xd, \
2336 p_yd, p_zd, p_volo, p_delx_xi, p_delx_eta, p_delx_zeta, p_delv_xi, p_delv_eta, \
2337 p_delv_zeta)
2338     for (i = 0; i < numElem; ++i) {
2339         const Real_t ptiny = 1.e-36;
2340         Real_t ax, ay, az;
2341         Real_t dxv, dyv, dzv;
2342
2343         const Index_t *elemToNode = &p_nodelist[8*i];
2344         Index_t n0 = elemToNode[0];
2345         Index_t n1 = elemToNode[1];
2346         Index_t n2 = elemToNode[2];
2347         Index_t n3 = elemToNode[3];
2348         Index_t n4 = elemToNode[4];
2349         Index_t n5 = elemToNode[5];
2350         Index_t n6 = elemToNode[6];
2351         Index_t n7 = elemToNode[7];
2352
2353         Real_t x0 = p_x[n0];
    
```

## Aspen code

```

147 kernel CalcMonotonicQGradients {
148     execute [numElems]
149     {
150         loads [8 * indexWordSize] from nodelist
151         // Load and cache position and velocity.
152         loads/caching [8 * wordSize] from x
153         loads/caching [8 * wordSize] from y
154         loads/caching [8 * wordSize] from z
155
156         loads/caching [8 * wordSize] from xvel
157         loads/caching [8 * wordSize] from yvel
158         loads/caching [8 * wordSize] from zvel
159
160         loads [wordSize] from volo
161         loads [wordSize] from vnew
162         // dx, dy, etc.
163         flops [90] as dp, simd
164         // delvx delvx
165         flops [9 + 8 + 3 + 30 + 5] as dp, simd
166         stores [wordSize] to delv_xeta
167         // delxi delxi
168         flops [9 + 8 + 3 + 30 + 5] as dp, simd
169         stores [wordSize] to delx_xi
170         // delvj and delvj
171         flops [9 + 8 + 3 + 30 + 5] as dp, simd
172         stores [wordSize] to delv_eta
173     }
174 }
    
```



# View Aspen performance models as normal performance analysis output with Gprof

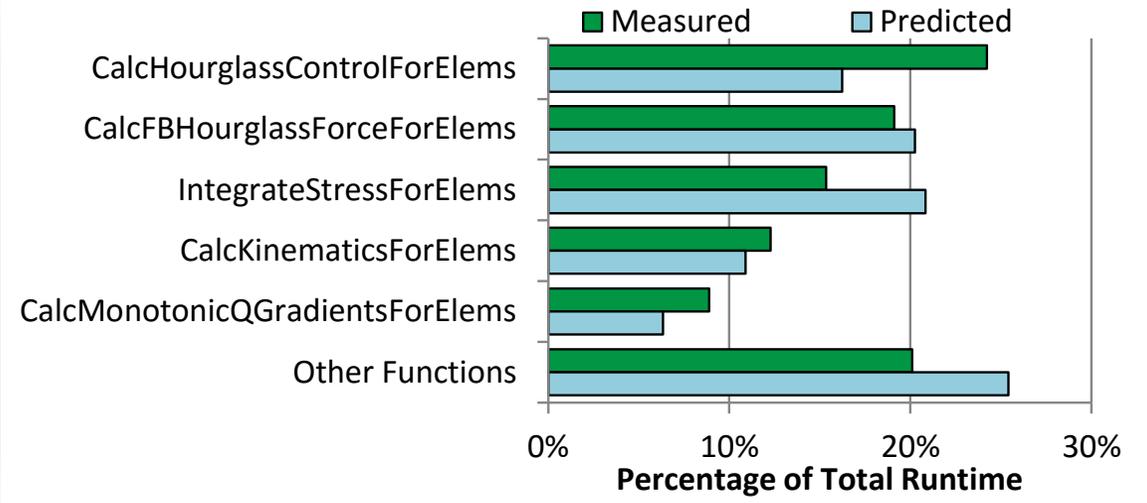
Flat profile:

% time	cum sec	self sec	calls	self ms/call	total ms/call	name
86.91	370.76	370.76	30	12358.52	12358.52	fft3d.localFFT
10.03	413.54	42.78	20	2139.09	2139.09	fft3d.exchange
3.06	426.57	13.03	20	651.71	651.71	fft3d.shuffle
0.00	426.57	0.00	10	0.03	0.03	exchange
0.00	426.57	0.00	10	0.03	0.03	buildNList
0.00	426.57	0.00	10	0.01	0.01	ljForce
0.00	426.57	0.00	30	0.00	0.00	integrate
0.00	426.57	0.00	10	0.00	42657.18	fft
0.00	426.57	0.00	10	0.00	42657.18	fft3d.main
0.00	426.57	0.00	1	0.00	426572.70	main

Not from gprof but rather aspen performance model

call graph:

index	%time	self	children	name
[ 1 ]	100.0	0.00	426.57	main [1]
		0.00	0.00	buildNList [8]
		0.00	0.00	exchange [7]
		0.00	426.57	fft [2]
		0.00	0.00	integrate [10]
		0.00	0.00	ljForce [9]
[ 2 ]	100.0	0.00	426.57	main [1]
		0.00	426.57	fft [2]
		0.00	426.57	fft3d.main [3]
[ 3 ]	100.0	0.00	426.57	fft [2]
		0.00	426.57	fft3d.main [3]
		42.78	0.00	fft3d.exchange [5]
		370.76	0.00	fft3d.localFFT [4]
		13.03	0.00	fft3d.shuffle [6]
[ 4 ]	86.9	0.00	426.57	fft3d.main [3]
		370.76	0.00	fft3d.localFFT [4]
[ 5 ]	10.0	0.00	426.57	fft3d.main [3]
		42.78	0.00	fft3d.exchange [5]
[ 6 ]	3.1	0.00	426.57	fft3d.main [3]
		13.03	0.00	fft3d.shuffle [6]
[ 7 ]	0.0	0.00	426.57	main [1]
		0.00	0.00	exchange [7]
[ 8 ]	0.0	0.00	426.57	main [1]
		0.00	0.00	buildNList [8]
[ 9 ]	0.0	0.00	426.57	main [1]
		0.00	0.00	ljForce [9]
[ 10 ]	0.0	0.00	426.57	main [1]
		0.00	0.00	integrate [10]



# Aspen Model User Queries

Benchmark	Runtime Order
BACKPROP	$H * O + H * I$
BFS	$nodes + edges$
CFD	$nelt * ndim$
CG	$nrow + ncol$
HOTSPOT	$sim_{time} * rows * cols$
JACOBI	$m_{size} * m_{size}$
KMEANS	$nAttr * nClusters$
LAPLACE2D	$n^2$
LUD	$matrix\_dim^3$
MATMUL	$N * M * P$
NW	$max\_cols^2$
SPMUL	$size + nonzero$
SRAD	$niter * rows * cols$

Table 2: Order analysis, showing Big O runtime for each benchmark in terms of its key parameters.

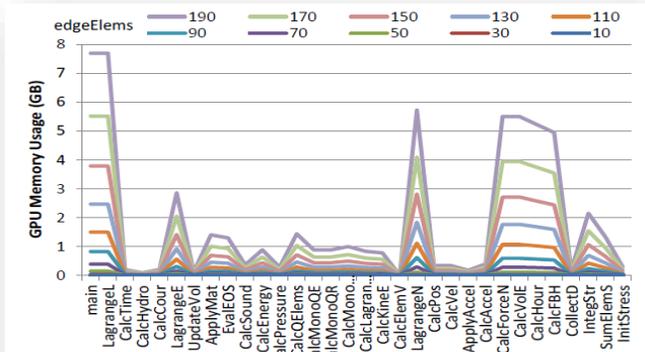
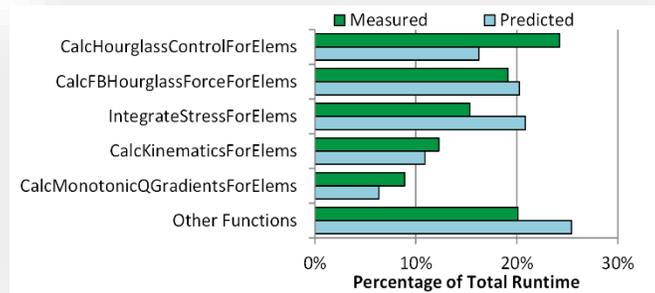


Fig. 8: GPU Memory Usage of each Function in LULESH, where the memory usage of a function is inclusive; value for a parent function includes data accessed by its child functions in the call graph.

Method Name	FLOPS/byte
InitStressTermsForElems	0.03
CalcElemShapeFunctionDerivatives	0.44
SumElemFaceNormal	0.50
CalcElemNodeNormals	0.15
SumElemStressesToNodeForces	0.06
IntegrateStressForElems	0.15
CollectDomainNodesToElemNodes	0.00
VoluDer	1.50
CalcElemVolumeDerivative	0.33
CalcElemFBHourglassForce	0.15
CalcFBHourglassForceForElems	0.17
CalcHourglassControlForElems	0.19
CalcVolumeForceForElems	0.18
CalcForceForNodes	0.18
CalcAccelerationForNodes	0.04
ApplyAccelerationBoundaryCond	0.00
CalcVelocityForNodes	0.13
CalcPositionForNodes	0.13
LagrangeNodal	0.18
AreaFace	10.25
CalcElemCharacteristicLength	0.44
CalcElemVelocityGradient	0.13
CalcKinematicsForElems	0.24
CalcLagrangeElements	0.24
CalcMonotonicQGradientsForElems	0.46

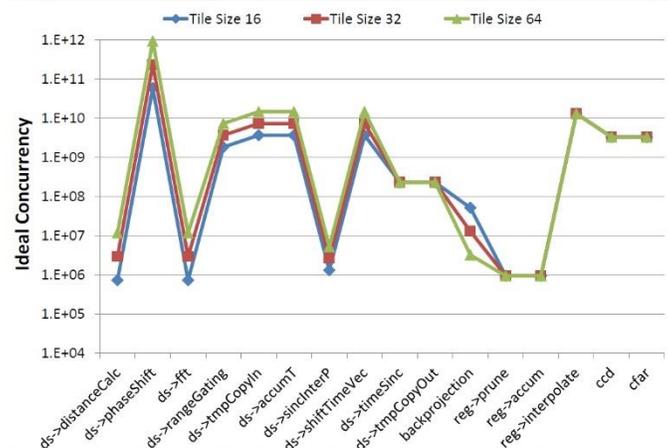


Figure 1: A plot of idealized concurrency by chronological phase in the digital spotlighting application model.

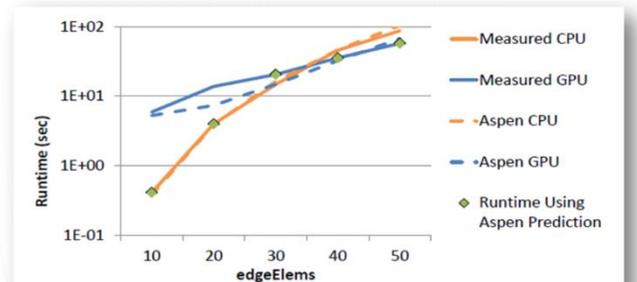


Fig. 7: Measured and predicted runtime of the entire LULESH program on CPU and GPU, including measured runtimes using the automatically predicted optimal target device at each size.

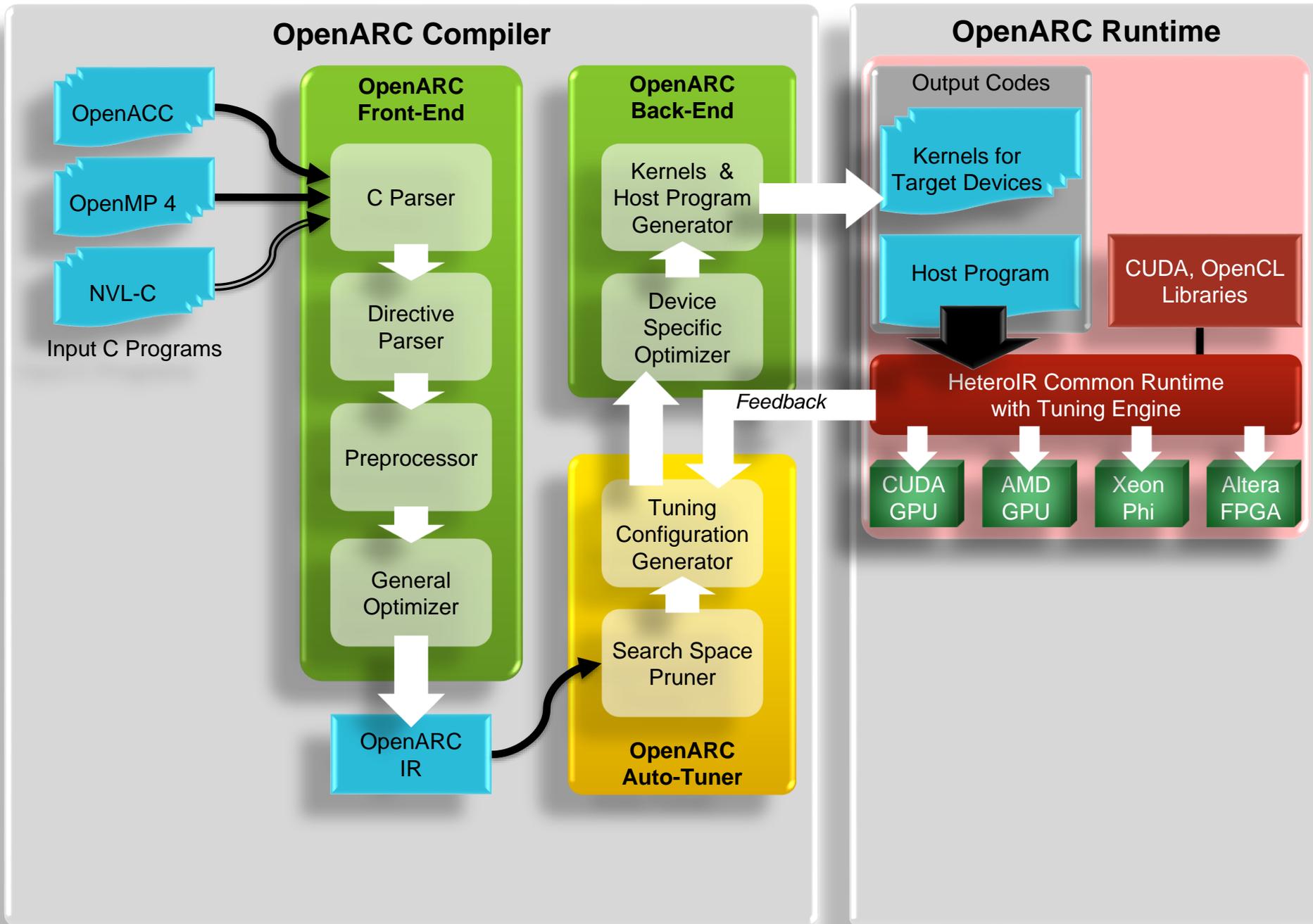
# Scheduling GPU Offloads with Aspen Performance Models

# Should the execution offload kernel to GPU or run on host CPU?

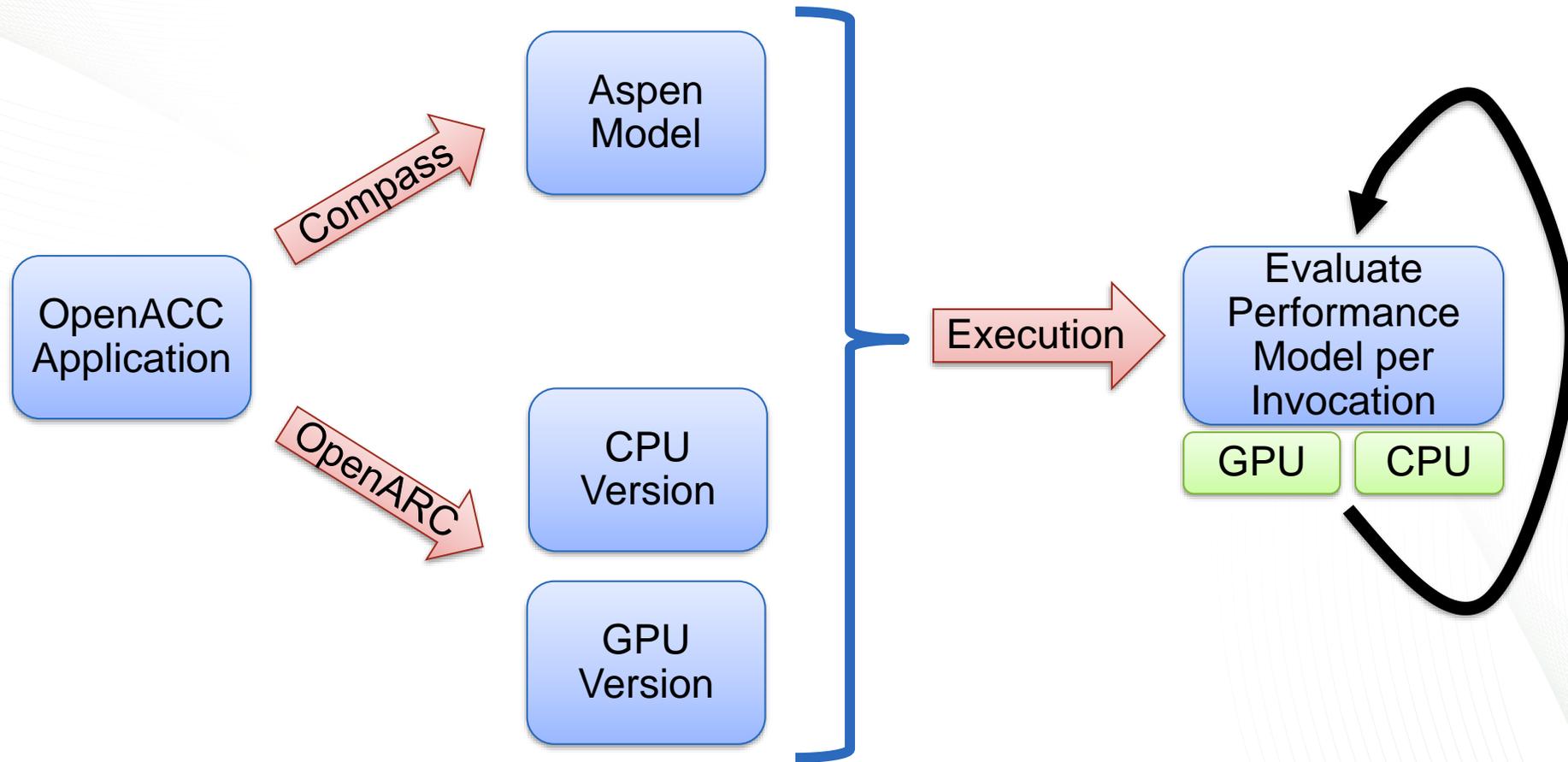
- Intuitively
  - When it is ‘small’, run the computation on the host CPU,
  - Otherwise, send it to the GPU
  - Expense of data movement (twice) and launch GPU kernels?
    - Depends on architecture
  - Simply offloading all computation is not smart
- Portability?
  - Need to account for performance, working set size, data transfer costs, ...

Listing 1: Input OpenACC Matrix Multiplication Code

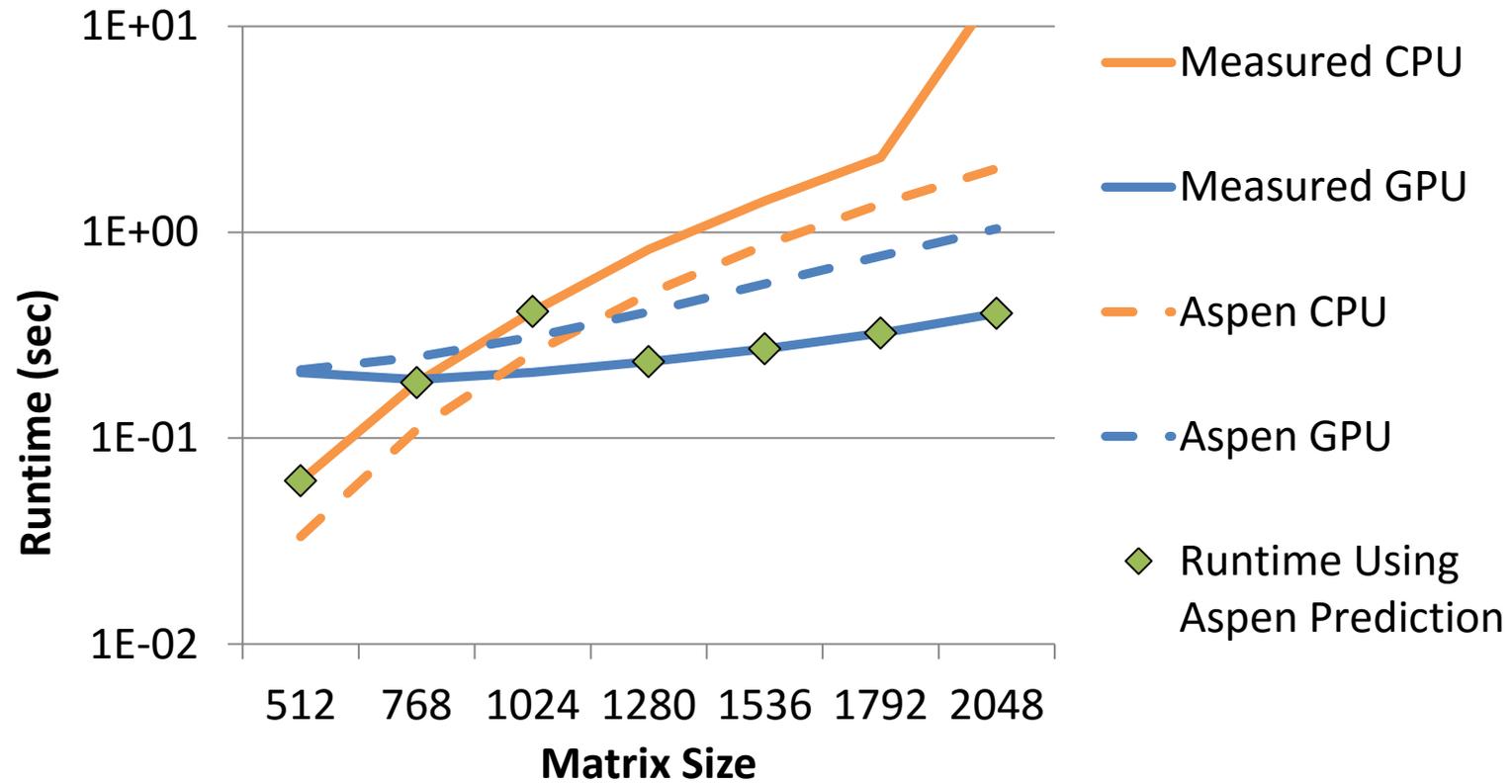
```
1 int N = 1024;
2 void matmul(float *a, float *b, float *c){ int i, j, k ;
3 #pragma acc kernels loop gang copyout(a[0:(N*N)]) \
4 copyin(b[0:(N*N)],c[0:(N*N)])
5   for (i=0; i<N; i++){
6     #pragma acc loop worker
7     for (j=0; j<N; j++) { float sum = 0.0 ;
8       for (k=0; k<N; k++) {sum+=b[i*N+k]*c[k*N+j];}
9       a[i*N+j] = sum; }
10  } //end of i loop
11 } //end of matmul()
12 int main() {
13   int i; float *A = (float*) malloc(N*N*sizeof(float));
14   float *B = (float*) malloc(N*N*sizeof(float));
15   float *C = (float*) malloc(N*N*sizeof(float));
16   for (i = 0; i < N*N; i++)
17     { A[i] = 0.0F; B[i] = (float) i; C[i] = 1.0F; }
18 #pragma aspen modelregion label(MM)
19   matmul(A,B,C);
20   free(A); free(B); free(C); return 0;
21 } //end of main()
```



# Informed Offloading Execution



# Matrix Multiply



# LULESH: Runtime and Working Set Size Predictions

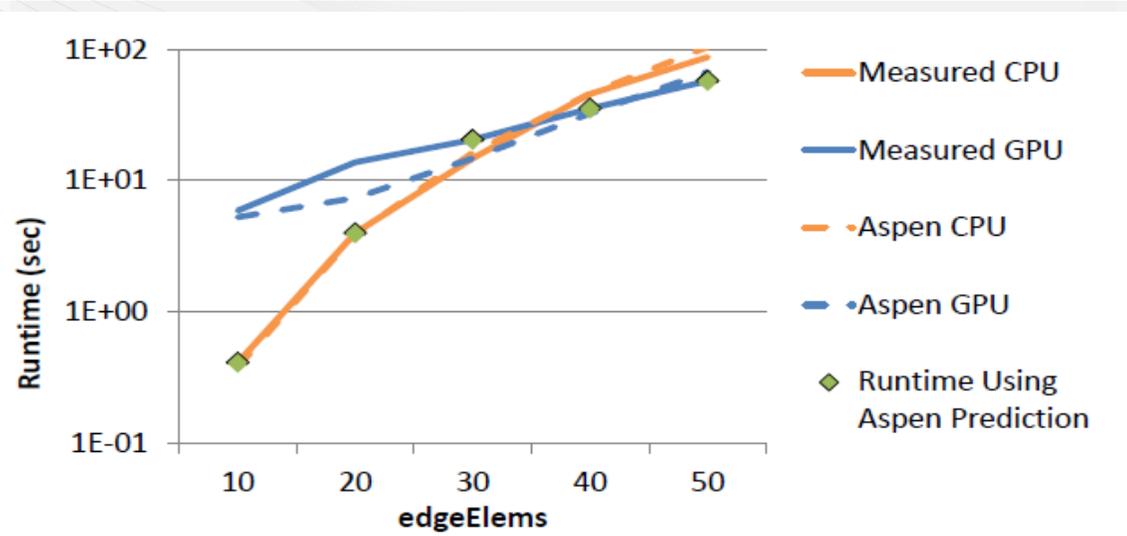


Fig. 7: Measured and predicted runtime of the entire LULESH program on CPU and GPU, including measured runtimes using the automatically predicted optimal target device at each size.

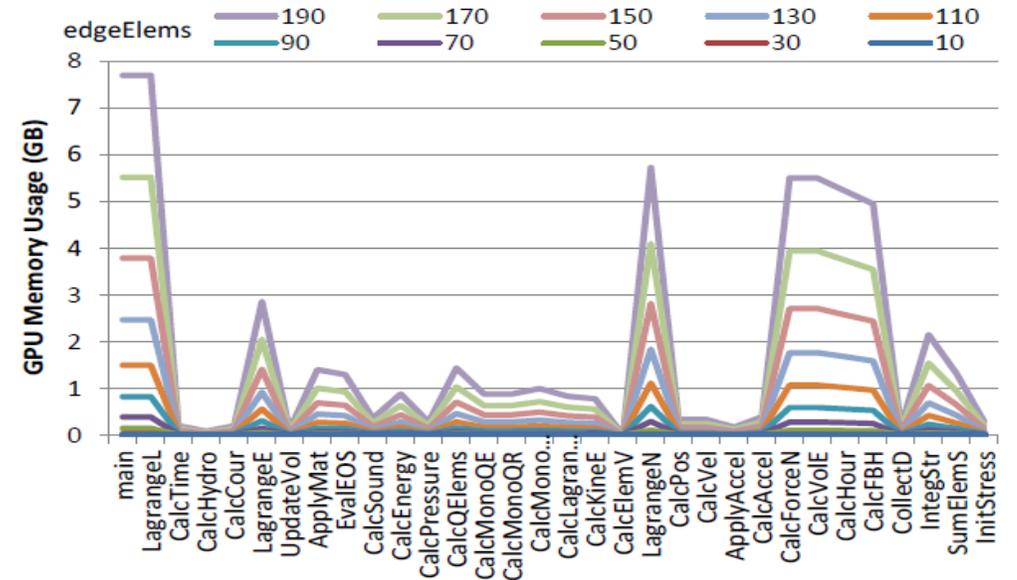
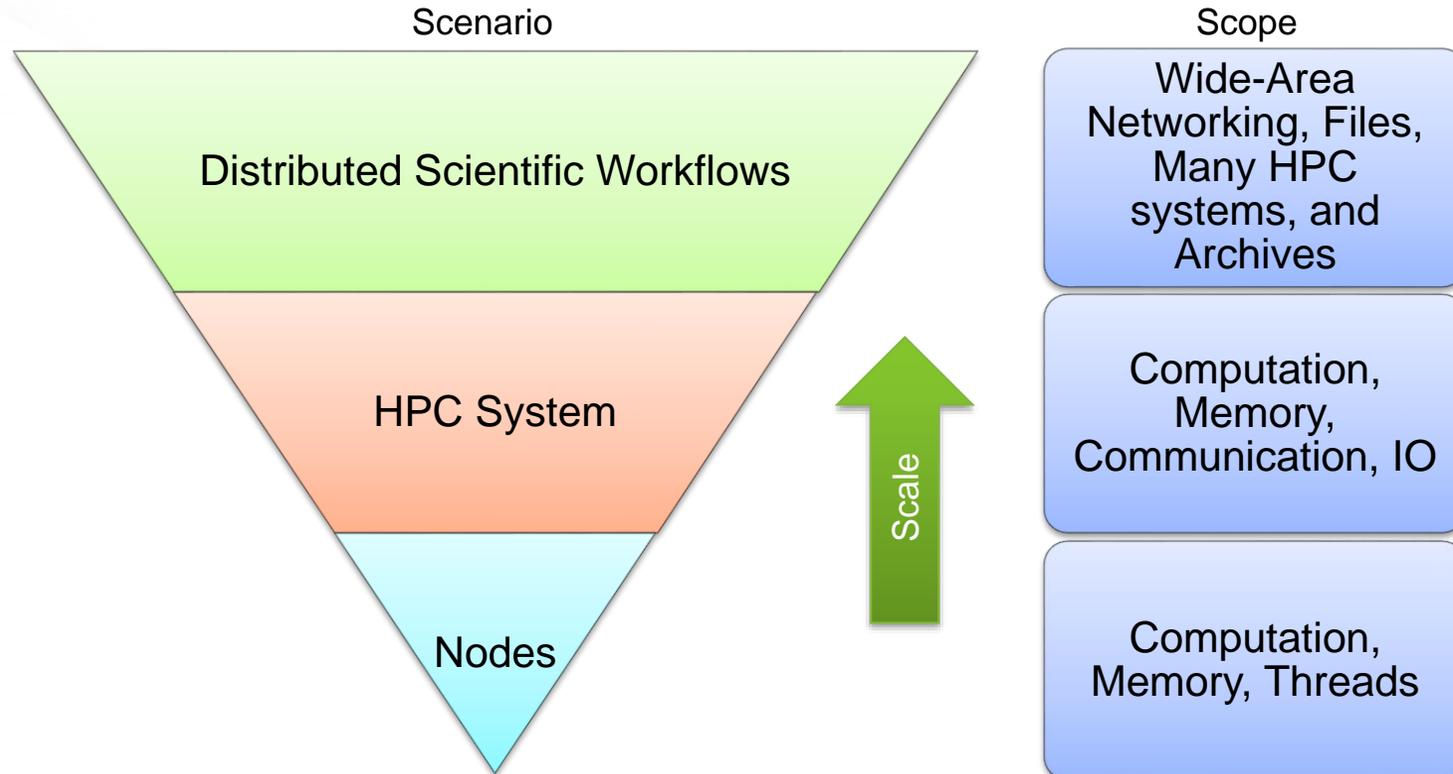


Fig. 8: GPU Memory Usage of each Function in LULESH, where the memory usage of a function is inclusive; value for a parent function includes data accessed by its child functions in the call graph.

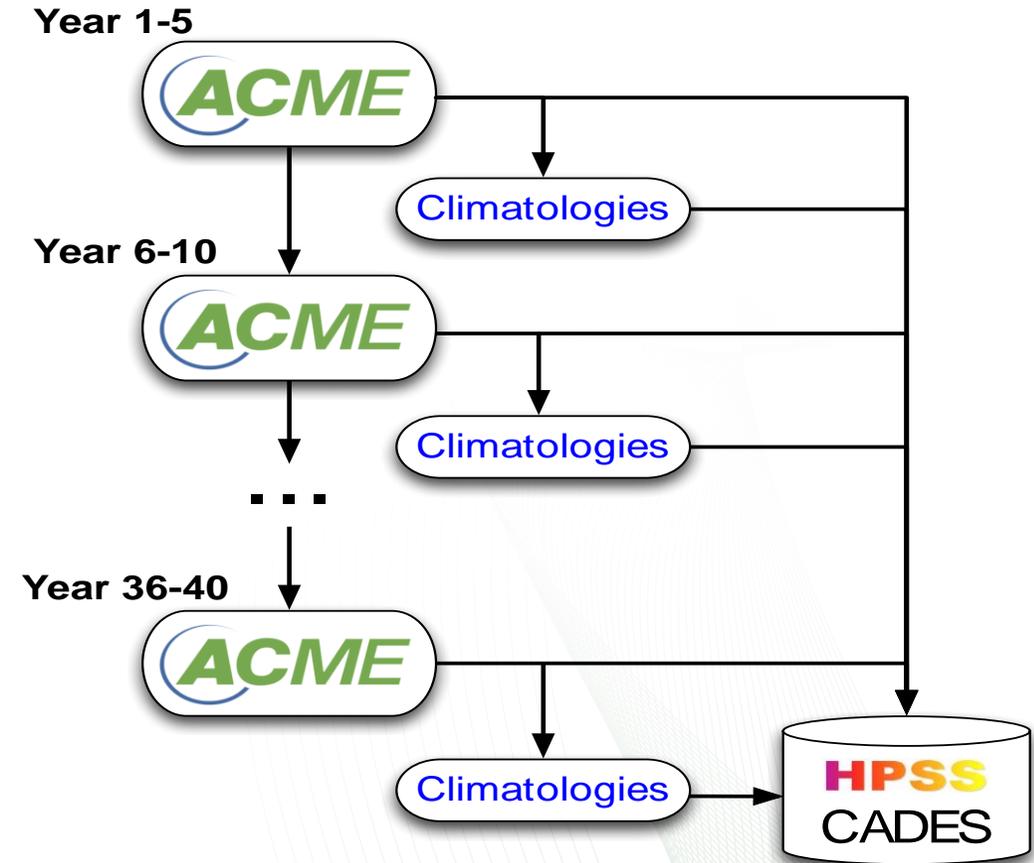
# Using Aspen for Distributed Workflows

# Aspen allows Multiresolution Modeling



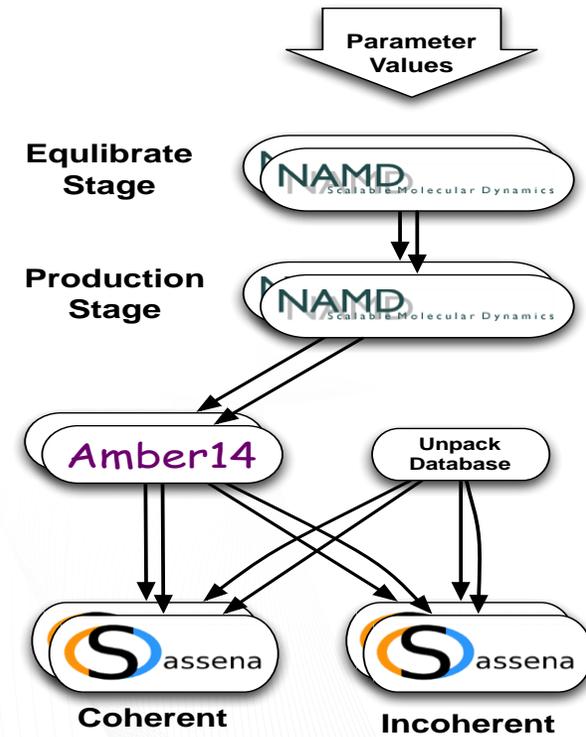
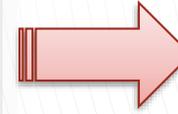
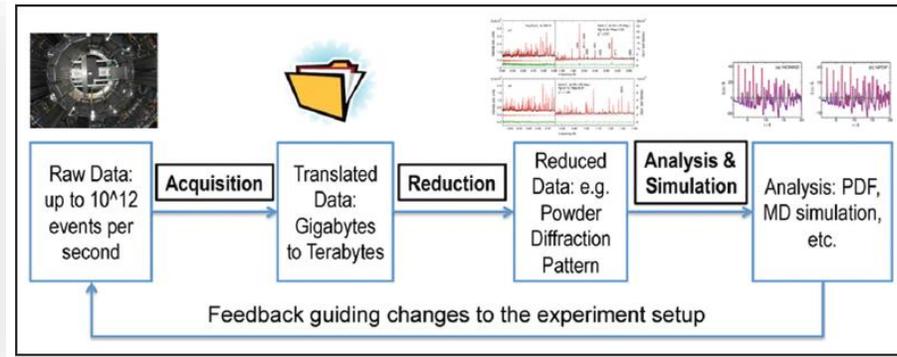
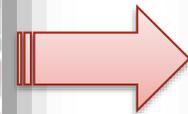
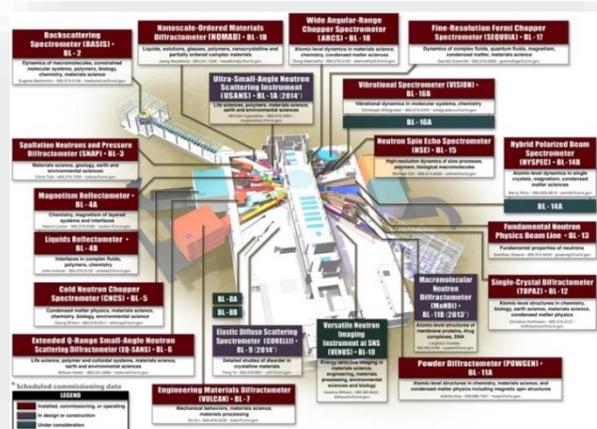
# Simulation: ACME Workflows

- Accelerated Climate Modeling for Energy (ACME)
- Coupled climate models with ocean, land, atmosphere and ice
- Climatologies and diagnostics give summaries of data
- Each stage of the workflow runs the ACME model for a few timesteps—helps keep simulations within batch queue limits
- Running on Hopper @ NERSC and Titan @ OLCF



# Experimental Data: SNS Workflows

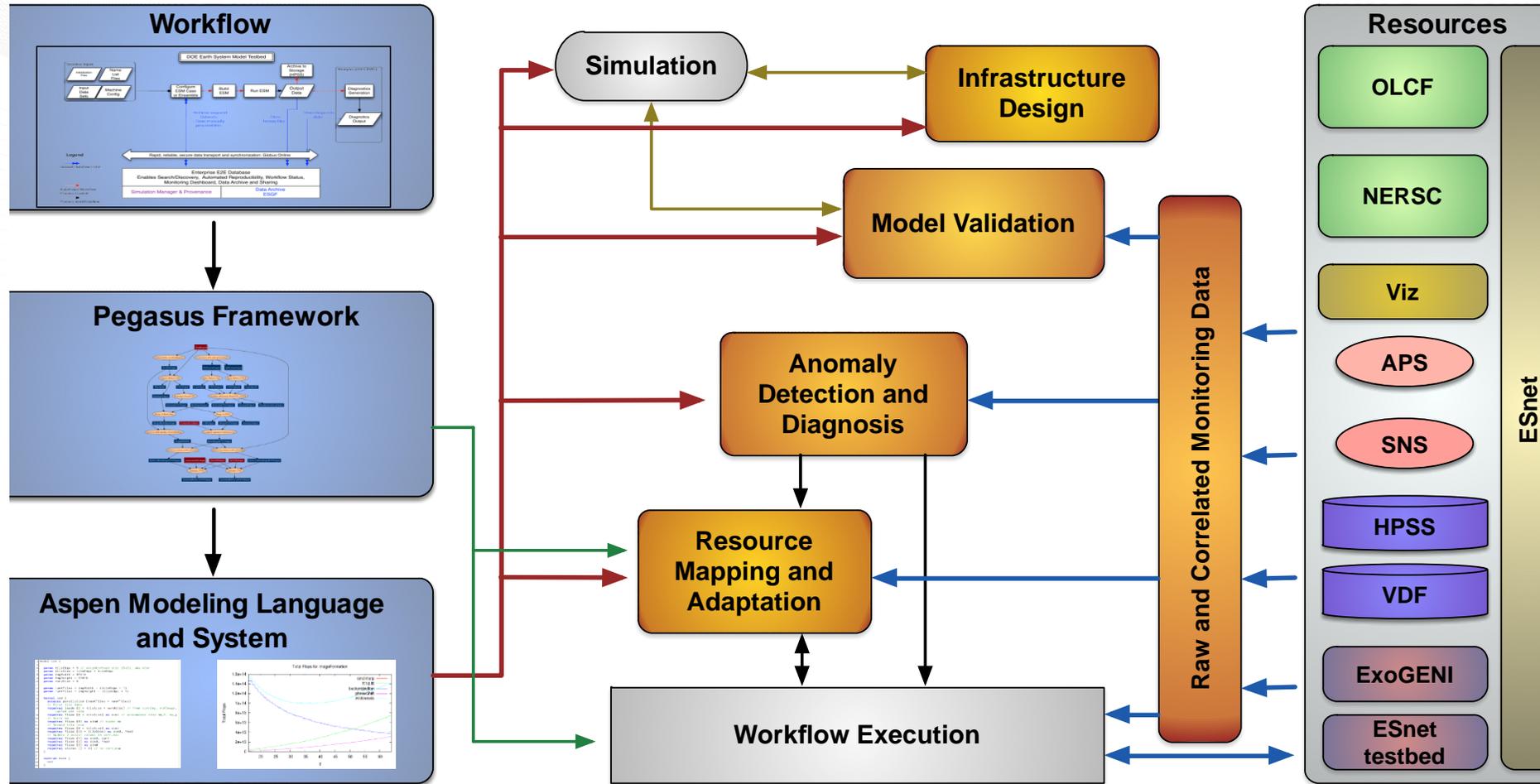
- Spallation Neutron Source
- Parameter sweep of molecular dynamics and neutron scattering
- Used to identify parameters that fit experimental data from SNS
- Currently being used for real science problems
- Large angle runs use 20 parameter values and require ~400,000 CPU hours
- Running on Hopper @ NERSC and coming to Titan @ ORNL



Post-processing and Viz



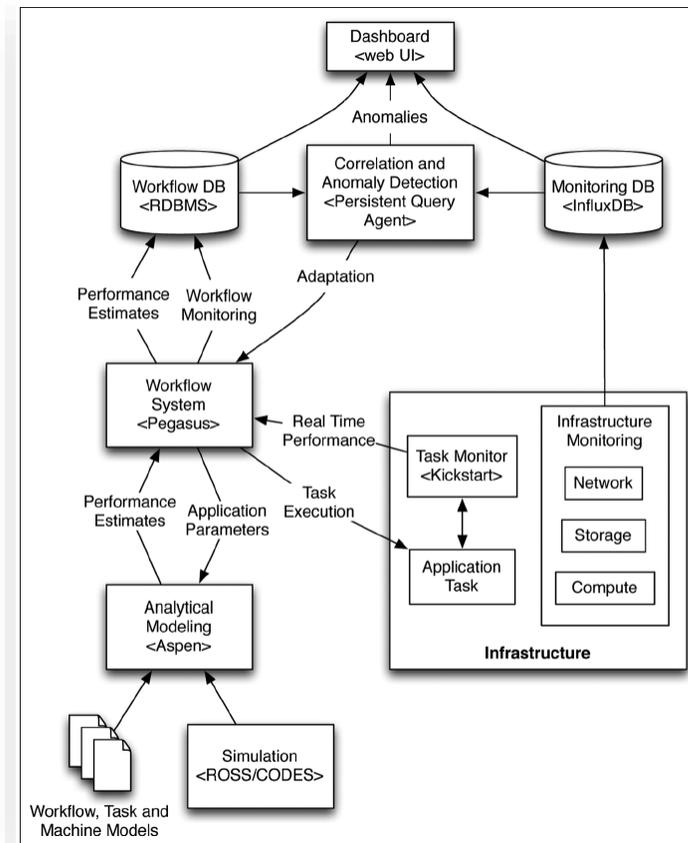
# PANORAMA Overview



# Automatically Generate Aspen from Pegasus DAX; Use Aspen Predictions to Inform/Monitor Decisions

```
1 kernel main
2 {
3   par {
4     seq {
5       call namd.eq_200()
6       call namd.prod_200()
7     }
8     seq {
9       call namd.eq_290()
10      call namd.prod_290()
11    }
12  }
13  par {
14    call unpack_database()
15    call ptraj_200()
16    call ptraj_290()
17  }
18  par {
19    call sassena_incoh_200()
20    call sassena_coh_200()
21    call sassena_incoh_290()
22    call sassena_coh_290()
23  }
24 }
```

Listing 1: Automatically generated Aspen model for sample SNS workflow.



# Executive Summary

- Both architectures and applications are growing more complex
  - Trends dictate that this will get worse, not better
  - This complexity creates irregularity in computation, communication, and data movement
- Dynamic resource management is one way to help manage this irregularity
  - Need accurate policies to guide resource decisions
  - Examples: greedy work stealing, algorithmic, historical, cost models, application specific, etc
- Posit that we can use application-specific performance models to inform scheduling decisions
  - Aspen performance modeling language helps create models
  - Two recent experiments
    - GPU offload
    - *Distributed scientific workflows*
- Q&A
  - *Can Aspen be accurate enough for these dynamic decisions?*
  - *How can we employ/influence the RT/OS in this process?*

# Acknowledgements



- Contributors and Sponsors

- Future Technologies Group: <http://ft.ornl.gov>
- US Department of Energy Office of Science
  - DOE Vancouver Project: <https://ft.ornl.gov/trac/vancouver>
  - DOE Blackcomb Project: <https://ft.ornl.gov/trac/blackcomb>
  - DOE ExMatEx Codesign Center: <http://codesign.lanl.gov>
  - DOE Cesar Codesign Center: <http://cesar.mcs.anl.gov/>
  - DOE Exascale Efforts:  
<http://science.energy.gov/ascr/research/computer-science/>
- Scalable Heterogeneous Computing Benchmark team:  
<http://bit.ly/shocmarx>
- US National Science Foundation Keeneland Project:  
<http://keeneland.gatech.edu>
- US DARPA
- NVIDIA CUDA Center of Excellence



# PMES Workshop @ SC16

- <https://j.mp/pmes2016>
- @SC16
- Position papers due June 17



Search this site

## 2016 Post-Moore's Era Supercomputing (PMES) Workshop Home

Co-located with [SC16](#) in Salt Lake City  
Monday, 14 November 2016

Workshop URL: <http://j.mp/pmes2016>  
CFP URL: <http://j.mp/pmes2016cfp>  
Submission URL (EasyChair): <http://j.mp/pmes2016submissions>  
Submission questions: [pmes16@easychair.org](mailto:pmes16@easychair.org)

**230**  
days until  
PMES Workshop @  
SC16

**News**  
[PMES Workshop Confirmed for SC16!](#)  
[Submissions open for PMES Position Papers on April 17](#)

**Important Dates**

- Submission Site Opens: 17 April 2016
- Submission Deadline: 17 June 2016
- Notification Deadline: 17 August 2016
- Workshop: 14 November 2016

This interdisciplinary workshop is organized to explore the scientific issues, challenges, and opportunities for supercomputing beyond the scaling limits of Moore's Law, with the ultimate goal of keeping supercomputing at the forefront of computing technologies beyond the physical and conceptual limits of current systems. Continuing progress of supercomputing beyond the scaling limits of Moore's Law is likely to require a comprehensive re-thinking of technologies, ranging from innovative materials and devices, circuits, system architectures, programming systems, system software, and applications.

The workshop is designed to foster interdisciplinary dialog across the necessary spectrum of stakeholders: applications, algorithms, software, and hardware. Motivating workshop questions will include the following. "What technologies might prevail in the Post Moore's

In cooperation with IEEE Computer Society  